



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁷ :

H04Q 11/00

A2

(11) International Publication Number:

WO 00/13456

(43) International Publication Date:

9 March 2000 (09.03.00)

(21) International Application Number: PCT/US99/20056

(22) International Filing Date: 31 August 1999 (31.08.99)

(30) Priority Data:

60/098,488	31 August 1998 (31.08.98)	US
60/118,874	5 February 1999 (05.02.99)	US
60/119,789	11 February 1999 (11.02.99)	US

(71) Applicant (for all designated States except US): ADERO, INC.
[US/US]; 432 Columbia Street, Suite 11, Cambridge, MA 02141-1039 (US).

(72) Inventors; and

(75) Inventors/Applicants (for US only): CROSBIE, David, B. [GB/US]; 47 Bay State Avenue, Somerville, MA 02144 (US). BAI, Joseph, J. [US/US]; 81 Meredith Road, Tewksbury, MA 01876 (US). McCORMICK, Kevin [US/US]; 253 Commonwealth Avenue, Boston, MA 02116 (US). GRAHAM, Chet, J. [US/US]; 12 White Horse Road, Hingham, MA 02043 (US). JONES, Ellis, Oliver [US/US]; 124 Rattlesnake Hill Road, Andover, MA 01810 (US).

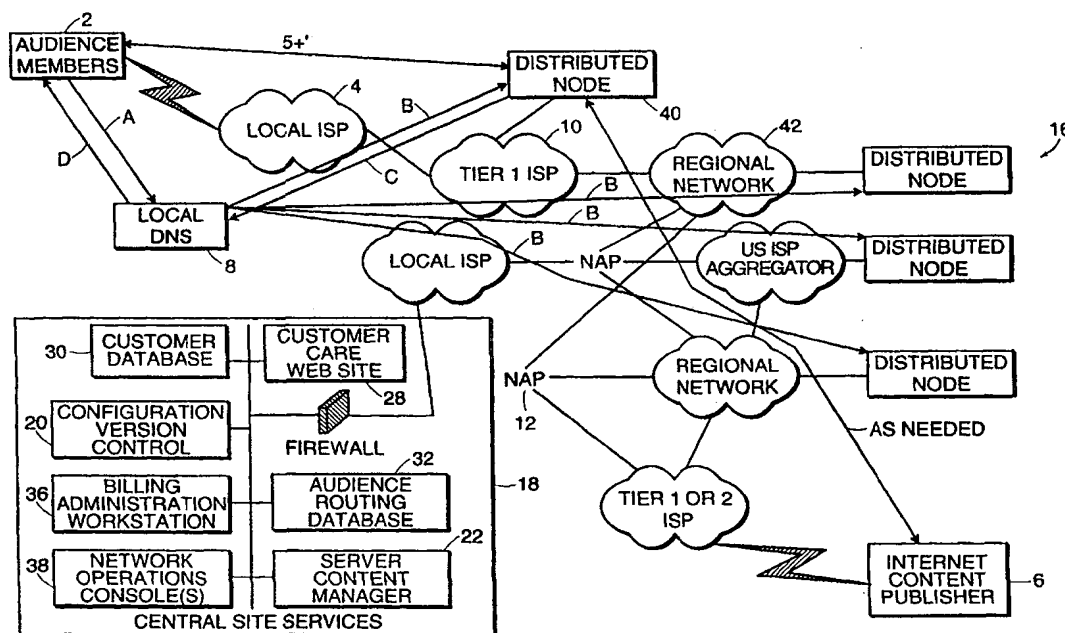
(74) Agents: WAKIMURA, Mary Lou et al.; Hamilton, Brook, Smith & Reynolds, P.C., Two Militia Drive, Lexington, MA 02421 (US).

(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

Published

Without international search report and to be republished upon receipt of that report.

(54) Title: METHOD AND APPARATUS FOR LOAD MANAGEMENT ON A COMPUTER NETWORK



(57) Abstract

The present invention relates to a global computer network system which includes distributed services and centralized services. The distributed services include a geo-traffic manager, a rebroadcaster and an audience data collector to optimally route data through the network from a content provider to an end user. The centralized services support the functionality of the distributed services.

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

METHOD AND APPARATUS FOR LOAD MANAGEMENT ON A COMPUTER NETWORK

RELATED APPLICATION(S)

This application claims the benefit of U.S. Provisional Application Nos.
5 60/119,789 filed on February 11, 1999, 60/118,874 filed February 5, 1999 and
60/098,488 filed on August 31, 1998, the entire teachings of which are incorporated
herein by reference.

BACKGROUND OF THE INVENTION

A multitude of users are accessing information stored on web servers located
10 anywhere in the world using computer networks such as the World Wide Web.
Referring to Fig. 1, an audience member or an end user 2 who is a person or a
computer system, typically uses a browser or a local Internet Service Provider (ISP)
4 to request and gain access to content from an Internet Content Provider or
Publisher (ICP) 6. In addition, a local Domain Name System (DNS) 8 server is used
15 by the audience member 2 as the name resolution protocol to convert host names to
Internet Protocol (IP) addresses. The Internet backbone in a particular geographic
region includes a collection of ISPs, and Tier 1 ISPs 10 which are large ISPs with
necessary resources to peer with a number of other large ISPs. Peering points are
points of interconnection between two autonomous IP networks such as ISPs. A
20 public Network Access Point (NAP) 12 is a place where ISPs of varying scale and
scope interconnect. A US ISP aggregator 14 is a

private peering service, within the US, used by the ICPs to maximize network reach
while avoiding the congestion of the public NAPs 12.

The world outside North America is experiencing an explosion in Internet usage. The problem is that accessing an overseas server can be a difficult and slow process because the International links are often more convoluted and congested than domestic ones. This has not been a key concern for most commercial web sites until recently because most content providers have aimed to satisfy their domestic market before expanding overseas. However, the global market is now as large as the North American one, and it is growing faster.

This demographic shift is too important for online content providers, most of whom are U.S. based, to ignore. An increasing number of Internet, Intranet, and Extranet content providers have growing populations of users in international markets where inadequate Internet infrastructure severely limits the quality of service. Though investment in international Internet infrastructure continues, it will not be able to keep up with the growth in demand.

The rapid growth of the Internet globally has given rise to severe performance problems for audience members who are remote with respect to Internet Content Publishers or Providers (ICPs). Performance is defined as the elapsed time between an audience member's request and the successful fulfillment of that request. The user's on-line experience is critical in ensuring that they stay on the site, return to the site, and complete transactions. The absolute speed of access is less important than the variability because web sites are designed with pages that download in a particular amount of time under normal conditions. It is therefore important to content providers that their site is consistently fast and reliable throughout their market.

Within the US market, market pressure has forced several novel solutions in three areas: backbone bandwidth, public peering, and access bandwidth congestion. Entire industries have arisen to provide additional backbone bandwidth, alternatives to public peering points, and burst capacity services. As a result, Internet performance within the US is currently a non-issue. Additionally, the US market benefits from being the leader in the development of Internet technologies and services to the point where most new offerings are tailored to US requirements. This is not the case outside the US. In the vast majority of markets, bandwidth remains the overwhelming constraint. Peering is either non-existent or highly restrictive and,

in many cases, takes place at congested US public peering points. Access bandwidth costs remain excessive. Additionally, these markets lag in the introduction of new Internet products and services.

Today, U.S. based ICPs may duplicate or mirror their own Web sites at
5 overseas locations, but this involves considerable effort. A mirror site is effectively a copy of the content and hardware (web servers, etc.) of the original site. Few Web hosting companies have a presence in every market. So in each target country, an ICP must find a suitable hosting firm or Internet Hosting Services (IHSs); negotiate a contract; build, install, and run the system; and pay a foreign currency invoice - all
10 the while dealing with language and cultural issues. In addition, the ICP faces the non-trivial challenge of keeping content replicated and synchronized, and centrally logging traffic.

Even the major step of building a mirror site only brings the content one step nearer to certain users. For example, a U.S. content provider setting up a mirror site
15 in Great Britain may not help French or Italian users at all.

Building a mirror site in every target market is not practical, especially if the demand is low or yet to be tested. There is very little difference in cost or effort between building a trial mirror site and a full installation.

Thus, outside the US, the challenges of doing business on the Internet make
20 it very difficult, if not impossible, for ICPs to get their message across to audience members due to performance issues. Interactions across geopolitical boundaries exacerbate this performance problem. In many cases, Internet traffic must flow through the US in order to pass between adjoining countries or in extreme cases within a country. Often the audience member finds it impossible, if not unbearable,
25 to view ICP content or to execute transactions with remote servers.

Historically, ICPs enter new markets by either ignoring the audience performance issues or by taking the extreme step of setting up a mirror site. As described hereinbefore, site mirroring is a costly proposition that tends to improve performance in one network-local market.

30 Recently, the first tier Internet Service Providers (ISPs) for example, providers such as AT&T, Sprint, GTE, etc., and Internet Hosting Services (IHSs) have discovered the global market. The ISP offerings fall short on market

penetration due to their competitive stature with respect to the local service providers, local with respect to a particular market. As a perceived threat, the ISPs find it difficult to establish the requisite peering relationships in many markets. The IHS entries fall short due to prohibitive real estate and bandwidth costs as well as the operational complexities associated with maintaining facilities in a large number of localities.

Another solution to this problem is caching, however while caching improves the user's experience it creates very serious problems for the content provider such as stale content and the elimination of information about the user. For example, some product companies produce commercial caching devices that can be placed closer to the user and so improve their on-line experience. But they do not have the software needed to control the content stored (i.e. removing unwanted pages) or retrieving the log files (files that record user activity). Some software companies provide some of this additional functionality. However, a content provider could buy caches from a product company and software from a software company but would still need to obtain server space and bandwidth in all their target markets from either an ISP (Internet Service Provider) in each location or as a package from an International ISP. They would then need to arrange maintenance and support in each country plus the 24 hour a day / 7 day a week ability to monitor and co-ordinate all the servers and the in-country support structure.

The resultant gap in service for the global market has given rise to opportunity for content distribution services. There are several incomplete solutions. For example, some companies have a service that out-sources the management of "flash floods" i.e. unusually high demand for service. Other companies have a similar service but only move the simple elements (i.e. just the images and single files) of the standard web page out to remote machines, while others place caches at the edge of their private global network but suffer from a very high underlying cost of service provision and limited penetration. This is due to utilizing the private network which is much more expensive than the Internet, hence this solution is cost prohibitive to most users. Therefore, there is still a challenge and need to minimize infrastructure costs for content distribution services while maximizing proximity of content to the target audience.

SUMMARY OF THE INVENTION

The present invention relates to global Internet performance and the ability to provide users with consistently fast and reliable access to a content provider's data.

A further object of the present invention is the ability to provide Internet user
5 monitoring and auditing companies with data on Internet usage. Another objective of the present invention is the ability to provide Internet network provision companies with the data they need in order to plan new network installations.

According to the present invention, a global system combines local content caching and mirroring with a traffic management mechanism. The system places
10 ICP content in close network proximity to their audience members and routes audience requests to the nearest content server, in terms of network topology. The traffic management system of the present invention in conjunction with the local content repositories is a solution to the issue of Internet audience performance.

The traffic management system is suited to the high-latency, high packet-loss
15 global Internet environment. Further, a geo-traffic manager is capable of scaling to the node count necessary to effectively reach the global Internet audience with acceptable performance. These two functions, local content caching and traffic management, necessarily form a tightly coupled system. One without the other is insufficient to fully address the performance issues of the global Internet.
20 Additionally, the highly distributed nature of the system provides greatly enhanced service availability. In addition to the elimination of single points of failure, the distributed node in this system of the present invention continues to operate correctly in the absence of communications with the central site.

In accordance with one aspect of the present invention, the system includes a
25 distributed node and central site services. The distributed node includes a geo-traffic manager, a rebroadcaster and an audience data collector. The distributed node is an extensible platform on which vertical applications can be built. The central site services include a network data collector, a customer care website, a billing administration workstation, a network monitor, a network operations console, a
30 customer database, an audience routing database, a configuration version control system and a server content manager.

In a preferred embodiment, the geo-traffic manager provides proximity routing of audience requests to rebroadcaster servers. The rebroadcasters provide the platform for local content distribution to global audience members. Additionally, the rebroadcaster supports content caching and mirroring based on configuration files constructed from data in the customer database. The rebroadcasters provide passive caching, pro-active caching and content mirroring.

BRIEF DESCRIPTION OF THE DRAWINGS

The foregoing and other objects, features and advantages of the invention will be apparent from the following more particular description of preferred embodiments of the invention, as illustrated in the accompanying drawings in which like reference characters refer to the same parts throughout the different views. The drawings are not necessarily to scale, emphasis instead being placed upon illustrating the principles of the invention.

Fig. 1 is an overview of a computer network environment in which the present invention solves problems of the prior art.

Fig. 2 is an overview of a computer network environment illustrating the system for load management in accordance with the present invention.

Fig. 3 is a block diagram illustrating the system interaction in accordance with the system for load management of the present invention.

Fig. 4 schematically illustrates domain name system traffic management used by the system in accordance with the present invention.

Figs. 5A and 5B illustrate the geo-traffic manager configuration and interfaces included in the system of the present invention.

Figs. 6A and 6B illustrate the rebroadcaster configuration and interfaces included in the system of the present invention.

Figs. 7A and 7B illustrate the audience data collector configuration and interfaces included in the system of the present invention.

Fig. 8 is a block diagram of the routing table configuration in accordance with the present invention.

Fig. 9 is a flow chart illustrating the rebroadcaster routing configuration in accordance with the present invention.

Fig. 10 is a diagram illustrating an exemplary use of the system in accordance with the present invention where an audience member requests data that is located within the rebroadcaster's cache.

Fig. 11 is a diagram illustrating another exemplary use of the system in accordance with the present invention where an audience member requests an
5 uncacheable content element.

Fig. 12 is a diagram illustrating another exemplary use of the system in accordance with the present invention where an audience member requests content not yet in cache.

10 Fig. 13 is a diagram illustrating another exemplary use of the system in accordance with the present invention where an audience member requests mirrored content.

DETAILED DESCRIPTION OF THE INVENTION

Referring to Fig.1, the system in accordance with the present invention is an
15 assemblage of two distinct service classes: distributed services 16 and central site services 18. A rebroadcaster, geo-traffic manager and audience data collector make up the distributed services 16 functions. A configuration control system 20, server content manager 22, network data collectors 24, network monitor 26, customer care web site 28, customer database 30, audience routing database 32, geographic data
20 table 34, a billing administration workstation 36 and a Network Operations Console (NOC) 38 make up the central site services.

In the present invention, distributed services nodes 40 are positioned at and serve as the Tier 1 ISP 10, Regional Network 42 and US ISP Aggregator 14. The Central Site Services 18 support the server functionality of the distributed services
25 16 resident in the distributed nodes 40.

The system in accordance with the present invention allows the start, expansion and end of service to a particular country simply by the flip of a software switch. Thus, the system can be used to test market demand without the delay, cost, and expense of a mirror site of the prior art. The system is cost effective for low
30 volumes of traffic as well as for very high volume but short duration events such as

trade shows, major company announcements and international sporting competitions. It can also act as a backup.

The system allows the content to be placed as close as possible to the end user while allowing the content provider to keep total control. The system is
5 effectively a distributed Web server network. It enables an online content provider's web server to be connected via the Internet, to a series of rebroadcast servers located in key countries.

Web pages consist of a series of elements, such as the menu bar, images, company logo, and text. Each element is transmitted from Web server to user's
10 browser as a separate file using a protocol called Hypertext Transfer Protocol (HTTP). The size of a typical page is growing as designers push the boundaries of Web design and faster Internet connections are introduced into homes and offices. Today a typical page consists of about 4-12 elements and totals around 40K bytes. Each page element can be divided into two types; static and dynamic.

15 Static content is the most striking part of a Web page. For example, it is the image of the product that the end user is buying, the navigation bar that shows where the user is in the buying process, and the button that the user presses to make the purchase. Even XML (eXtensible Markup Language), the new generation of HTML (Hyper-Text Markup Language), contains significant amounts of static content.

20 Dynamic content is normally the information that is provided uniquely for a user. It is the part of the page that shows how many items are in the user's electronic shopping basket, and how much the user has spent so far. In contrast to the static content, the dynamic content files are very small, and normally represent less than 20% of the total amount of data needed to make up the page. Most dynamic content
25 is derived from back end systems such as a stock inventory database or an airline reservation system. These are very difficult and expensive items to duplicate not only because of the hardware and software required, but also because of the difficulty of synchronizing the different versions of the same database.

Typically 80% or more of the data that makes up a Web page is static. The
30 system of the present invention positions the content closer to the end user by keeping the static content at the rebroadcast site and providing the dynamic content via an optimized Internet link.

The optimized link is a way of transmitting data reliably and efficiently over the Internet even when the route from audience member or end user to ICP is convoluted and a significant proportion of the transmitted packets of data are lost.

The typical Internet infrastructure consists of a series of communication links
5 which are interconnected by routers. Data is sent as packets and each packet is forwarded from one router to the next until the packet reaches its destination or is lost. The router decides on the appropriate path by consulting a database called a routing table. This table essentially says "if you receive a packet destined for these end locations then send it to this particular router next". The problem with this form
10 of routing is that it is essentially static and often setup to minimize Internet backbone operator's costs rather than optimizing performance.

Applicants have discovered that by combining information from the Internet routing tables with information on the location of Internet congestion, it is possible to relay the customer's data around the congestion. In effect, the customer uses the
15 private network included in the system of the present invention which in turn uses the Internet as its network. In some circumstances where it is not possible to consistently route around a particular problem, dedicated or reserved communications links are required to bridge the problem area.

Fig. 2 illustrates the system for load management on a computer network in
20 accordance with the present invention. An audience member or user 2 makes a request via a browser (local Internet Service Provider (ISP)) 4 to a local DNS server 8. The audience member 2 selects a universal resource locator (URL), and a DNS lookup is initiated with the local DNS server 8. If the local DNS server does not hold a valid record for the target server name (selected URL), the local DNS
25 transmits a datagram request to all authoritative distributed nodes 40 for the name. The nearest distributed node 40 having a name server responds to the local DNS 8 first and then determines which rebroadcaster is available. The distributed node 40 returns to the local DNS service 8 a response in the form of an IP address identifying the rebroadcaster that can service the request. The local DNS (LDNS) caches the
30 response according to the time-to-live (TTL) associated with the identified rebroadcaster name (IP address). The TTL field indicates how long the IP address should be cached by the local DNS before being discarded. TTL also gives the DNS

server an indication of how recent and therefore how reliable the information they receive about other hosts is. The local DNS 8 ignores all subsequent responses from other servers. The local DNS forwards the received and cached rebroadcaster's IP address to the audience member 2. The audience member 2 then initiates
5 communication directly with the rebroadcaster housed in a distributed node 40. It should be noted that although Fig. 2 illustrates the scenario of the geo-traffic manager in the distributed node that returns a response in the form of an IP address is the same distributed node housing the available rebroadcaster, a response could have originated from another distributed node. A standard hypertext transfer
10 protocol (HTTP) is used to make a request to the rebroadcaster. If necessary, the rebroadcaster contacts the ICP's origin server to retrieve un-cached content elements. It should be noted that the ICP's/customers cede responsibility for the name subdomain that is hosted on the rebroadcasters.

The method of the present invention is advantageous to an end user as it
15 allows the user to directly communicate with a rebroadcaster once the IP address of an available rebroadcaster is provided to the end user by a geo-traffic manager. The rebroadcaster then communicates with the ICP to provide the desired content to the end user. A portion of the content, such as the static content, is moved into the rebroadcaster which improves the speed and reliability for the user. Further, the
20 method of the present invention is advantageous to the ICP as the ICPs do not have to build expensive mirror sites.

Fig. 3 illustrates the system interaction between the distributed services 16 or nodes 40 and the central site services 18. Each distributed node includes a respective rebroadcaster 50, geo-traffic manager 52 and audience data collector 54
25 which are each discussed later. Although the node 40 support all three functions, the system is capable of operating with the functions physically separated from one another. Any grouping of the functions within a distributed node is, in fact, supportable by the system due to the network transparent nature of the system. The functions of the Internet Content Publisher 6, Local DNS 8 and audience member 2
30 are outside the scope of control of the distributed nodes 40 of the present invention.

Continuing with Fig. 3, the Configuration Version Control System 20 serves as the central point of control and audit for server configuration files throughout the

system in accordance with the present invention. Used in conjunction with the Server Content Manager 22, the Configuration Version Control System 20 facilitates full roll-forward and rollback functionality. Rollforward functionality is the configuration version control needed for releasing new versions of the configuration to production systems. Rollback functionality is the configuration control needed for referencing a previous version released to production systems. The version control system simplifies change tracking, auditability and rollback in the event of a misconfiguration. The configuration version control system receives data from the Customer Database 30 and the audience routing database 32 (discussed later).

10 The Network Data Collectors 24 (Figs. 1 and 3) serve as an independent set of eyes into the global Internet. These collectors gather data on network connectivity, latency and routing, DNS response times, as well as HTTP performance. The Configuration Version Control system 20 oversees the configuration files for the collectors. The configurations are the result of queries against the Customer Database 30 indicating: tests to be performed, regions to be tested, networks to be tested, and URLs to be probed.

 The Customer Database 30 acts as the repository for the test results. Test results contribute to customer reporting, and geo-traffic manager routing tables.

20 The Network Monitor 26 is the point of concentration for all system monitoring logging with the exception of the Network Data Collector 24. The operating system platform of the Network Data Collectors is incompatible with the server agents that forward data to the monitor. In another particular embodiment, the operating system platform can be compatible with the server agents that forward data to the monitor.

25 The Network Operations Console 38 provides a high-level summary-view of the present invention system. It also gives the operators of the invention system an interface to control system components through the server agents and the Customer Database 30. The Network Operations Console 38 makes changes to the Customer Database 30 to reflect availability of the rebroadcaster 50. The changes are reflected in the geo-traffic manager routing tables 64.

30 The Customer Care Web Site 28 allows the ICP 6 to alter a number of parameters that influence the provisioning of the present invention service. The

customer (ICP) may make changes in near real-time to a number of service parameters, including:

- URLs served
- Regions, countries, and Rebroadcasters purchased
- 5 • Premiums services selected
- Reporting
- Network data collection
- Billing
- Account characteristics

10 Customer/ICP initiated changes are limited programmatically to those allowed contractually and appropriate to the current service level. In addition to the customer interfaces within the Customer Care Web Site 28, additional levels of access exist for Sales, Operations, Administration, and Engineering personnel. These interfaces support such functions as customer addition, password reset, and
15 content access control definition.

The Customer Database 30 houses customer contact, account, service configuration, and usage data with respect to the end users 2 using the customer or ICPs 6 content. Usage by the end user 2 can be tracked by measures such as pages served from local cache and requests routed to the ICP's origin server. The value of
20 a transaction between an ICP 6 and a user 2 can be monitored or calculated and the Customer/ICP can be charged accordingly. The Customer Database 30 is the source of the rebroadcaster configuration data as well as a source of data for the geo-traffic manager's content routing tables. This database contains information on customer URL caching and mirroring on a per rebroadcaster basis. It also acts as the
25 repository for each geo-traffic manager 52, rebroadcaster 50, and Network Data Collector 24 logs. In a preferred embodiment, this log data arrives at the Customer Database 30 over a database link from a WebSpective Manager database serving as the Network Monitor 26.

In a preferred embodiment, the geo-traffic manager 52, rebroadcaster 50 and
30 audience data collector 54 that are built from the Customer Database 30 are checked into a Configuration Version Control system 20 prior to distribution to the target servers by WebSpective's binary coded decimal (BCD) service. The version control

system simplifies change tracking, auditability and rollback in the event of a misconfiguration. Additionally, all changes to the Customer Database are logged in an audit log.

In a preferred embodiment, the customer database server is based on the Sun
 5 PCI SPARC platform from the Concorde Group (2x366mHz, 1GB RAM, 2x9GB
 internal and 6x18GB hot-swappable disks). It makes use of secured version of
 Solaris 2.6. Disk partition management is overseen by Sun Solstice Disk Suite
 (SDS). SDS provides striping and mirroring capabilities. WebSpective's v3.0 agent
 collects system statistics to the central WebSpective Manager database (Oracle 8I).
 10 WebSpective's BCD service transports configuration files for the services on the
 rebroadcaster. Netscape FastTrack 3.0.1 acts as the web server for administrative
 access.

The audience routing database 32 serves as the repository for all network
 topology data. A number of routing feeds contribute to the database including the
 15 Internet Routing Data Feed 60, the Geographic Data 34 and the Network Exception
 Data 62. Several classification parameters help in analyzing the relative weight of
 the particular datum in defining the topology portion of the geo-traffic manager
 routing table. In particular, the relative weight of a data point is set as a function of
 volatility and age or time. With respect to time weighting an exponentially decaying
 20 weight to data points is applied. As for volatility, weighting can be achieved by at
 least two different approaches to obtain the proper smoothing. The first approach is
 time domain averaging which is used to avoid skewing the calculation. The

equation used is
$$\bar{x} = \frac{\sum_{i=1}^n x_i}{\Delta t} \text{ (for } x < x_{\max} \text{)}.$$

The second approach used is a frequency domain averaging in which an
 25 exponentially decaying weight is applied to data points below a specific threshold.

Given the topology of the network, the existence of a geographic data table
 may seem counterintuitive. Local regulations and a desire to target content drive the
 need to incorporate such data into audience routing decisions.

The Server Content Manager 22 acts as the central control point for data distribution throughout the invention system with exception of the Network Data Collector 24. The operating system platform of the Network Data Collectors is incompatible with the server agents that accept data from the Server Content Manager 22. The Server Content Manager 22 reports on distribution status to the Network Monitor 26 and in turn to the Network Operations Console 38.

In a particular embodiment, the Billing Administration Workstation 36 does not play an active role in the load management, data distribution and trafficking operation of the present invention but instead provides business automation between the owner/operators of the system and the Customers/ICPs. The Billing Workstation takes summary usage data from the Customer Database 30 to generate summary and detailed statements. In a preferred embodiment, the Billing Workstation is based on NT and Visual Basic, and communicates with the Customer Database 30 through ODBC over Transmission Control Protocol/Internet Protocol (TCP/IP).

Illustrated in Fig. 4 is the Domain Name System (DNS) based traffic management in accordance with the system of the present invention between an audience member 2 and a subject distributed node 40. The geo-traffic manager 52 provides intelligent proximity routing of audience 2 requests to rebroadcasters 50. The geo-traffic manager 52 performs the audience routing function by answering domain name system (DNS) requests for name resolution from the audience member's local DNS server 8. A User Datagram Protocol (UDP) is used to communicate between the audience member and the geo-traffic manager. The IP address of the optimum rebroadcaster is sent to the audience member by the geo-traffic manager 52. The audience member 2 then communicates directly with the rebroadcaster 50 to access the desired content.

DNS is the most efficient means of performing this function because it adds no additional overhead in the interaction between the audience member and the target content. In order to utilize this DNS resolution scheme a Customer/ICP must delegate authority for the sub-domain global geo-traffic managers.

Referring to Fig. 5A, strategically placed geo-traffic managers act as the authoritative DNS server for publishers' sub-domains. The geo-traffic manager 52 includes the function blocks of content mirroring and configuration files 66; a geo-

traffic manager server 68, an HTTP server 70, a server agent 72 for data distribution, collection and management, an operating system 74, Redundant Array of Independent Disks (RAID) 76 which is a method of storing same data in different places on multiple disks to improve input/output operating, and server hardware 78.

- 5 The Server Content Manager 22 controls transportation of the configuration files and mirrored content to the Server Agents on the geo-traffic manager system. The geo-traffic manager server interfaces with the audience LDNS and the audience data collector 54.

Referring to Fig. 5B, in a particular embodiment, the public BIND server 77
10 interfaces with the audience/user LDNS. The public BIND server 79 accepts the LDNS request from the audience and performs a look-up for the ICP origin server having the target content and the available rebroadcaster. The public BIND server returns the IP address of the optimum rebroadcaster to the LDNS. The rebroadcaster Bind Server provides a centrally administered means of resolving origin server IP
15 address for rebroadcaster's. This service is important for the retrieval of cache content by the rebroadcasters and for routing requests for uncacheable content.

Multiple geo-traffic managers act as authoritative DNS servers for the desired sub-domain in a 'first response' DNS server configuration. The 'first response' configuration provides optimal responsiveness to the LDNS, as well as a
20 high level of redundancy and therefore improved availability. In Fig. 5B, horizontal distance does not correspond to either network or geographic distance. In a particular embodiment, the DNS server uses BIND 8.2 with modifications to query a proprietary key-value data structure. This data structure incorporates a number of factors:

- 25
- Internet Routing Data (in the form of Classless Internet Domain Routing(CIDR) tables)
 - Network performance data (from the customer database and the network data collector)
 - Geographic Data
 - 30 • Rebroadcaster deployment data and Customer configuration records (from the customer database)

- Server availability and load (from the customer database, network monitor, and the network operations console)
- The out-of-band audience data (from the customer database and audience data collector).

5 The out-of-band mechanism provides audience specific network data to fine tune the geo-traffic manager routing tables. This mechanism has several advantages over existing, polled-agent products as it provides fast local lookup, the local lookup provided is insensitive to link latency and loss, and the local lookup provided is highly scaleable.

10 Data from the Network Data Collectors 24, the Audience Data Collector 54, the Network Monitor 26, and the Network Operations Console 38 reside in the Customer Database 30 with the customer and server configuration data. The distillation process combines the data from the Customer Database 30 with data from the Internet Routing Data Feed 60, Geographic Data 34, and a Network
15 Exception Table 62. The distillation process checks the resultant table structures into a version control system prior to distribution to the distributed nodes 40 by the Server Content Manager 22.

 This mechanism is relatively insensitive to the LDNS honoring time-to-live (TTL) on the DNS lookups for proximity routing. Since subsequent audience
20 requests, to the same LDNS, are highly likely to yield the same target rebroadcaster 50, local DNS caching of the response rarely leads to stale data caching. Local DNS caching becomes an issue when a rebroadcaster 50 becomes unavailable due to an equipment or network fault. Here a DNS that does not honor a small TTL will continue to send audience members 2 to a rebroadcaster 50 that has been marked as
25 down within the geo-redirector's tables. As all modern versions of DNS servers support TTL functionality, only servers specifically configured to ignore the parameter will exhibit the aberrant behavior. It is possible to detect the signature of poorly configured LDNS server through geo-traffic manager log analysis.

 In the distributed node 40 configuration where the geo-traffic manager 52
30 and rebroadcaster 50 are co-resident, each distributed node acts as a DNS server. Using first response in this configuration provides first tier routing for the LDNS. Once the distributed node receives the DNS request, the table lookups can place less

weight on network proximity, relying on network latency during the look up to determine proximity. In a split configuration, routing table proximity plays a more active role in determining the best rebroadcaster 50 to service the request. Of course, in any configuration the relative weighting of the factors in a routing
5 decision is fully configurable.

In a preferred embodiment, the geo-traffic manager 52 is based on the Sun PCI SPARC platform from the Concorde Group (1x366MHz, 512MB RAM, 2x9GB internal disks). It makes use of secured version of Solaris 2.6. Disk partition management, striping, and mirroring is overseen by Sun SDS. WebSpective's v3.0
10 agent collects system statistics and DNS server logs to the central WebSpective Manager database. WebSpective's BCD service transports configuration files for the DNS services on the geo-traffic manager. Netscape FastTrack 3.0.1. acts as the web server for administrative access. The system in accordance with the present invention uses modified bind answers DNS request from local DNS servers and
15 from rebroadcasters.

The geo-traffic manager 52 illustrated for carrying out proximity routing of audience requests to rebroadcasters 50 of the present invention is purely exemplary. Other traffic managers can be utilized in light of the teachings herein.

Illustrated in Fig. 6A is the configuration of the rebroadcaster 50. The
20 rebroadcaster provides data collection, local load management, content distribution, content cache, content mirroring, and web servers. The rebroadcasters 50 are distributed throughout the network/Internet to provide localized content services globally.

The rebroadcaster 50 provides the platform for local content distribution to
25 global audience members 2. The rebroadcaster includes the function blocks of content mirroring and configuration files 80; an HTTP cache server 82, an HTTP server 84, a server agent 86 for data distribution, collection and management, an operating system 88, RAID 90 and server hardware 92. The rebroadcaster 50 supports content caching and mirroring based on configuration files constructed
30 from data in the Customer Database 30. The Configuration Version Control server 20 acts as a point of control and audit for the rebroadcaster configuration files. The

Server Content Manager 22 controls transportation of the configuration files and mirrored content to the Server Agents on the rebroadcaster system.

Referring to Fig. 6B, a commercially available HTTP caching technology forms the basis for the caching services on the rebroadcaster 50. The Server Agent 86 collects HTTP Cache Server 82 log information and transports the data to the Network Monitor 26 database (and subsequently to the Customer Database 30). The agent 86 also offers an administrative interface for cache invalidation and pre-population. In an exemplary embodiment, cache pre-population is a three-step process:

1. Spider the ICP's site to build a complete content tree listing.
2. Distribute the list to the relevant rebroadcasters through the Server Content Manager 22 and the Server Agent 86.
3. Execute a job on each of the rebroadcasters to request each of the URLs that make up the site.

In a preferred embodiment, the rebroadcaster is based on the Sun PCI SPARC platform from the Concorde Group (2x366MHz, 512MB RAM, 2x9GB internal and 3x18GB hot-swappable disks). It makes use of secured version of Solaris 2.6. Disk partition management is overseen by Sun SDS. SDS provides striping and mirroring capabilities. WebSpective's v3.0 agent collects system statistics to the central WebSpective Manager database (Oracle8I). WebSpective's binary coded decimal (BCD) service transports configuration files for the services on the rebroadcaster. Netscape FastTrack 3.0.1 acts as the web server for mirrored content and for administrative access. Inktomi's Traffic Server 3.0 supports content caching.

In a preferred embodiment, caching services on the rebroadcaster 50 are based on Inktomi's Traffic Server 3.0. WebSpective's BCD supplies configuration files from the central CVS production configuration repository. WebSpective's agent collects Traffic Server log information and transports the data to the central database. WebSpective's agent also offers an administrative interface for cache invalidation and pre-population.

The Inktomi traffic server 3.0 illustrated for carrying out the caching function of the present invention is purely exemplary. Other forward caching hardware or software can be utilized in light of the teachings herein.

Mirroring involves copying an entire website including all the logging,
5 access control, and databases to computers in a new geographic location and then putting in place the systems needed to keep the databases synchronized. For content mirroring services, in a preferred embodiment, WebSpective's BCD transports client content to the rebroadcasters from a central BCD server located, for example, at HarvardNet. This content is stored locally in the UNIX file systems and is served by
10 a Netscape FastTrack 3.0.1. webserver. The WebSpective server agent monitors and manages the Netscape Server. It also collects the web server logs to the central WebSpective Manager database.

ICP content delivery to the server content manager is performed by a number of mechanisms: File Transfer Protocol (FTP) sent to an authenticated server, SCP
15 and HTTP.

Referring to Figs. 7A and 7B, the configuration and interfaces of the audience data collector are illustrated. The audience data collector 54 provides a data collection service along with the Network Data Collector 24. The audience data collector includes the function blocks of content mirroring and configuration files
20 100, an audience data collector server 102, an HTTP server 104, a server agent 106 for data distribution, collection and management, an operating system 108, RAID 110 and server hardware 112. The Server Content Manager 22 controls transportation of the configuration files and mirrored content to the Server Agents 106 on the audience data collector system.

25 In contrast to the Network Data Collectors 24 which collect data on the ICP origin server 6, the rebroadcasters 50, and the network connectivity between selected ISPs, the audience data collectors 54 collect data on network connectivity to actual Customer/ICP networks and ultimately Customer/ICP IP addresses, if desired. In another preferred embodiment, the audience data collector 54 can operate as a stand-
30 alone service on independent hardware as an alternative to the illustrated embodiment with the audience data collector integrated in the distributed node 40.

Th audience data collectors receive configuration data both from the Customer Database 30 and directly from the geo-traffic manager 52. The latter feeds near real-time routing and response time data for actual audience members 2. Data collected by the audience data collectors 54 feed into the central Customer Database 30 and is used for both routing table calculations and reporting services. This data is essential in the development of an accurate view of global network peering. Time average of the audience data collector data is a factor in the geo-traffic manager 52 configuration changes. This slow feedback avoids changes in response to short-lived network events. Similar to all the other distributed components 50, 52, the audience data collector 54 is managed by the server agent 106 and provides a web administrative interface for network operations.

Fig. 8 illustrates the routing table configuration in accordance with the system of the present invention. The calculation of the geo-traffic manager routing tables 64 (Fig. 3) aggregates a number of data sources into a format searched efficiently by the BIND process. Data from the Network Operations Console 38, the audience data collector 54, the Network Monitor 26, and the Network Data Collectors 24 resides in the Customer Database 30. The Router Table Distillation Process 120 combines the data from the Customer Database 30 with data from the Audience Routing Database 32 which includes data from the Internet Routing Data Feed 60, Geographic Data 34 and the Network Exception Data 62. This data then forms an input into the geo-traffic manager routing tables 64 (shown also in Fig. 3). The resultant table structures form an input into the Configuration Version Control System 20 prior to distribution to the geo-traffic managers 52 in the distributed nodes 40 by the Server Content Manager 22.

Factors contributing to the geo-traffic manager routing tables include:

- Rebroadcaster availability and load
- Audience Data Collector-based routing data
- Network Data Collector-based routing data
- Systems' configuration from the Network Operations Console
- Customer input data from the Customer Care Web Site
- System deployment and configuration data from the Customer Database

- Internet Routing Data
- Geographic Data
- Network Exception Data

The router table distillation process condenses this data into a key-value
5 structure that is suitable for use by the BIND function call.

Fig. 9 illustrates the rebroadcaster configuration in accordance with the system of the present invention. The calculation of the rebroadcaster configuration is simpler than the calculation of the geo-traffic manager routing tables 64. Rebroadcaster configurations 63 are the formatted results of queries against the
10 Customer Database 30. Data relevant to this process arrives primarily from the Customer Care Web Site 28. Thus, data from the Customer Care Web Site forms an input into the Customer Database 30. Queries against the Customer Database in step 122 from the Rebroadcaster Configuration files 63. These configurations form an input into the Configuration Version Control System 20 for change tracking and
15 auditability purposes. The configurations are distributed to the rebroadcasters 50 resident in the distributed nodes 40 by the Server Content Manager 22.

Fig. 10 illustrates a preferred embodiment showing a characteristic interaction between audience members 2 and the system of the present invention as well as interactions between the system of the present invention and the ICPs 6 web
20 site when required. The figure illustrates a use case where an audience member 2 requests data that is located within the rebroadcaster's 50 cache. Here the interaction is similar to that between an audience member and a local web server.

The audience member 2 makes a DNS request for a target content using a local DNS 8 which is forwarded to the geo-traffic manager 52. The geo-traffic
25 manager 52 after performing look-up for the ICP 6 having the target content and an optimal rebroadcaster 50, returns a DNS response to the LDNS in the form of an IP address of the available rebroadcaster which is forwarded to the audience member 2. The audience member directly communicates via HTTP with the available rebroadcaster identified. The rebroadcaster returns the target content located within
30 the rebroadcaster cache to satisfy the audience member request.

In Fig. 11, the audience member 2 requests an un-cacheable content element by making a DNS request similar to the process described in Fig. 10. Personalized

pages and stock quotes are examples of un-cacheable content. Here the rebroadcaster 50 transparently (to the audience member) forwards the request to the ICP/origin server 6, accepts the response, and forwards the requested data to the audience member 2. The system then includes the loading of local display elements
5 that make up the remainder of the page.

In Fig. 12, the audience member requests a previously un-cached content element by making a DNS request similar to the process described in Fig. 10. The rebroadcaster 50, once again, transparently forwards the request to the ICP/origin server 6, accepts the response, and forwards the requested data to the audience
10 member 2.

Fig. 13 illustrates an exemplary case where the content requested by the audience member 2 by making a DNS request similar to the process described in Fig. 10, is served from a local site mirror on the rebroadcaster 50. The dotted arrow
130 indicates the batch, bulk, content replication mechanism used to populate the
15 mirror. In this exemplary case, if the content is not on the rebroadcaster 50, the audience member 2 receives an error condition. With properly constructed content, this should never be the case since the mirror should behave identically to the ICP/origin server 6. A particular embodiment supports mixing of mirrored, cached, and uncacheable content through a technique of URL rewrite.

20 While this invention has been particularly shown and described with references to preferred embodiments thereof, it will be understood by those skilled in the art that various changes in form and details may be made therein without departing from the scope of the invention encompassed by the appended claims.

CLAIMS

What is claimed is:

1. In a computer network formed of a plurality of communication channels and a plurality of digital processors coupled to the communication channels for communication thereon, computer apparatus for traffic management
5 comprising:
a rebroadcaster operating on one of the digital processors to optimally route data through the network from a content provider on one of the digital processors to a user on one of the digital processors; and
10 a geo-traffic manager providing proximity routing of user requests to the rebroadcaster based on a plurality of criteria to provide optimal routing of data.
2. The computer apparatus of Claim 1 wherein the geo-traffic manager is located on the digital processor used by the rebroadcaster.
- 15 3. The computer apparatus of Claim 1 wherein the geo-traffic manager is located on a different digital processor than the rebroadcaster.
4. The computer apparatus of Claim 1 wherein the criteria further comprises network proximity to available sources of content.
5. The computer apparatus of Claim 4 wherein the network proximity criteria is
20 a function of where the rebroadcasters are located.
6. The computer apparatus of Claim 4 wherein the network proximity criteria is a function of how close a user is to the rebroadcaster.
7. The computer apparatus of Claim 4 wherein the network proximity criteria is a function of whether target content is available on the rebroadcaster.

8. The computer apparatus of Claim 1 wherein the criteria is a function of geography.
9. The computer apparatus of Claim 1 wherein the criteria is a function of rebroadcaster availability.
- 5 10. The computer apparatus of Claim 1 wherein the criteria is a function of rebroadcaster load.
11. The computer apparatus of Claim 1 wherein the criteria is a function of out of band data collection regarding network latency.
12. The computer apparatus of Claim 1 wherein the criteria is a function of
10 network link types.
13. The computer apparatus of Claim 1 further comprising central services systems.
14. The computer apparatus of Claim 13 wherein the central services system includes a configuration control system to control and audit server
15 configuration files.
15. The computer apparatus of Claim 13 wherein the central services system includes a server content manager to distribute data.
16. The computer apparatus of Claim 13 wherein the central services systems includes a plurality of network data collectors to collect data on network
20 parameters.
17. The computer apparatus of Claim 16 wherein the network parameters includes network connectivity, latency and routing.

18. The computer apparatus of Claim 16 wherein the network parameters includes DNS response times and performance of hypertext transfer protocol (HTTP).
- 5 19. The computer apparatus of Claim 13 wherein the central services systems includes a network monitor to monitor logging activity.
20. The computer apparatus of Claim 13 wherein the central services systems includes a customer database to store customer contact, account, service configuration and usage data.
- 10 21. The computer apparatus of Claim 13 wherein the central services systems includes a billing workstation.
22. The computer apparatus of Claim 13 wherein the central services systems includes a customer care website to allow the Internet Content Provider (ICP) to alter a plurality of parameters.
- 15 23. The computer apparatus of Claim 22 wherein the parameters include URLs served.
24. The computer apparatus of Claim 22 wherein the parameters include services selected.
25. The computer apparatus of Claim 13 wherein the central services systems includes an audience routing database to store network topology data.
- 20 26. The computer apparatus of Claim 13 wherein the central services systems includes a network operations console to provide an interface to control systems components.

27. The computer apparatus of Claim 1 wherein the rebroadcaster functions to provide data collection, local load management, content distribution, content caching, and content mirroring.
28. In a computer network formed of a plurality of digital processors coupled to communicate with each other, a method for managing traffic comprising the steps of:
- 5 a user processor selecting a Universal Resource Locator (URL);
initiating a domain name server (DNS) lookup with a local DNS (LDNS) server in response to the user processor selecting a URL;
10 the LDNS transmitting a request to a plurality of geo-traffic managers for name of a target server to the URL;
a geo-traffic manager responding to the LDNS server with the Internet Protocol (IP) address for an optimal rebroadcaster;
the LDNS server forwarding the rebroadcaster's IP address to the user processor; and
15 the user processor initiating communication with the rebroadcaster to request particular content on the internet.
29. The method of Claim 28, wherein the LDNS server caches the response regarding the IP address of the nearest rebroadcaster.
- 20 30. The method of Claim 28, wherein upon receiving the IP address of the nearest rebroadcaster, the LDNS ignores all subsequent responses from other geo-traffic manager.
31. The method of Claim 28, wherein the optimal rebroadcaster is chosen as a function of at least one of proximity, availability of content, geography,
25 availability of rebroadcaster, rebroadcaster load, and network latency.
32. In a computer network having a plurality of digital processors coupled to communicate with each other, at least one being an Internet Content Provider

(ICP), a method to provide fast and reliable communication between the ICP and other digital processors on the network comprising:

using a database to track the lowest latency route through the network between two or more digital processors in the network;

5 according to the database, using a rebroadcaster selected from a plurality of rebroadcasters as a relay to route data optimally through the network from a content provider to a digital processor on the network;

 determining the optimum rebroadcaster in the plurality of rebroadcasters using data gathered from monitoring flow of traffic combined
10 with network routing.

33. The method of Claim 32 further comprising producing digital processor usage profile by using data from the aggregate traffic flows through the network of rebroadcasters.

34. The method of Claim 32 further comprising tracking digital processor usage
15 by such measures as pages routed to the digital processor.

35. The method of Claim 32 further comprising calculating the value of a transaction between an ICP and a digital processor and charging the ICP according to the value.

36. The method of Claim 32 wherein the digital processor is an international
20 digital processor.

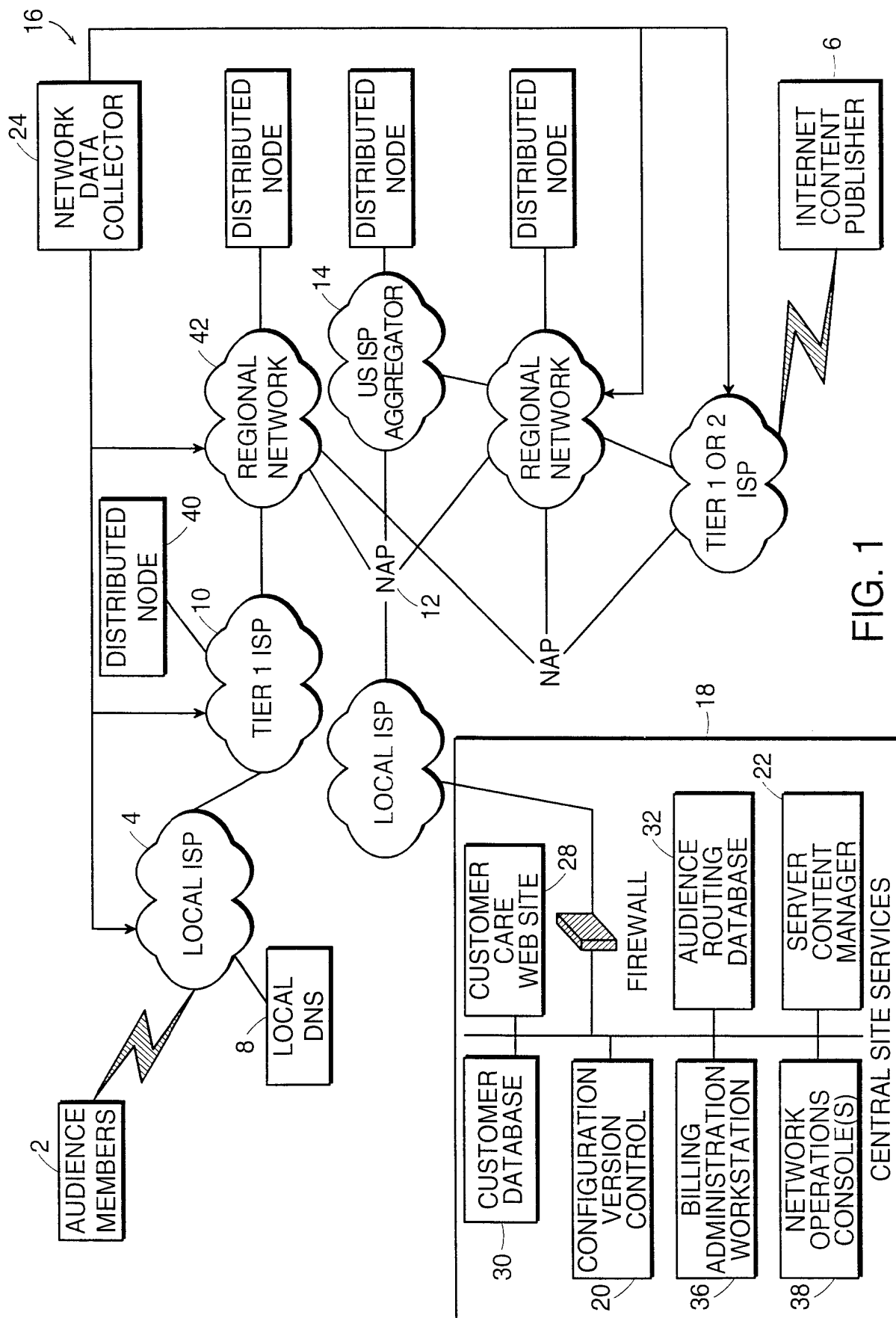
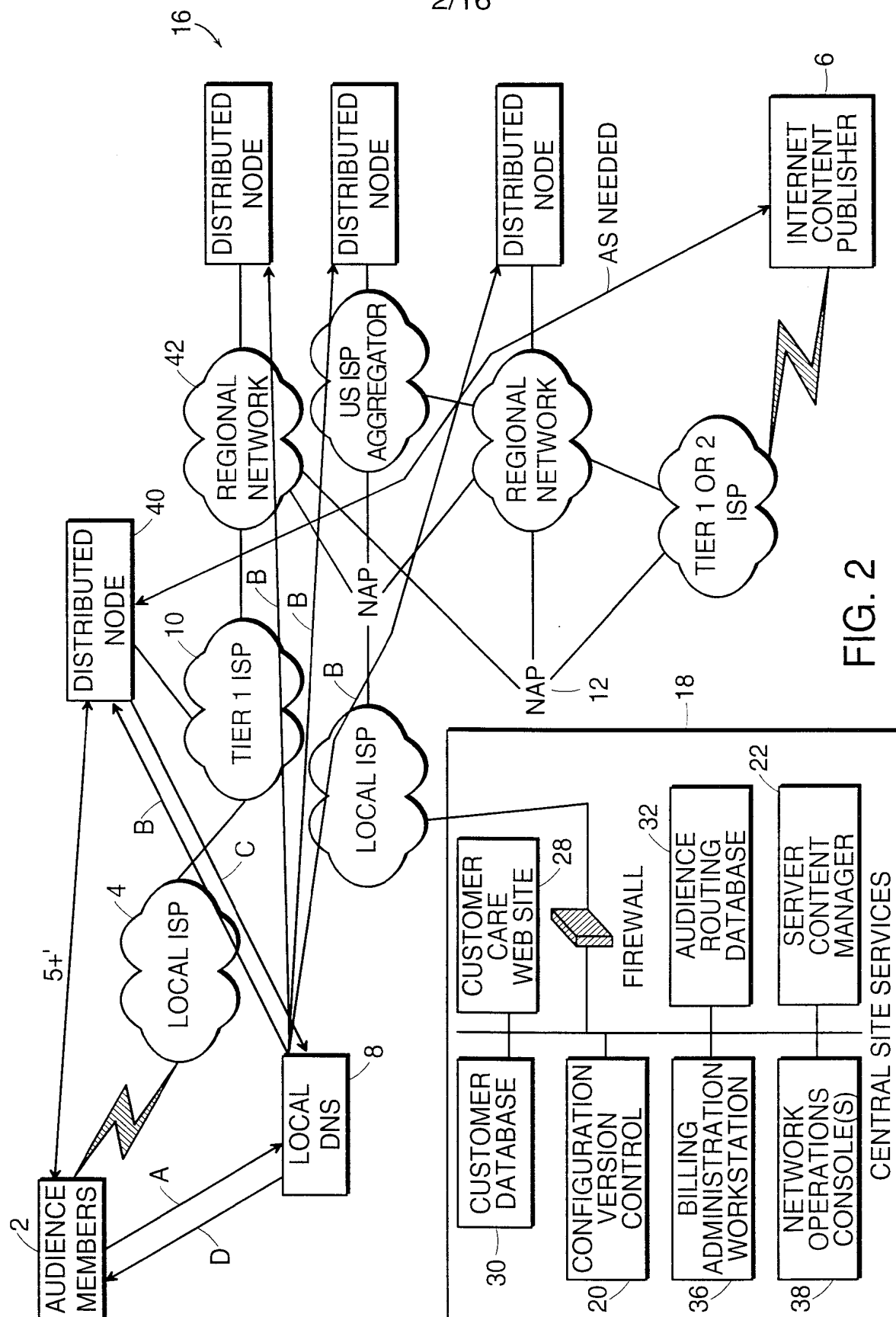


FIG. 1



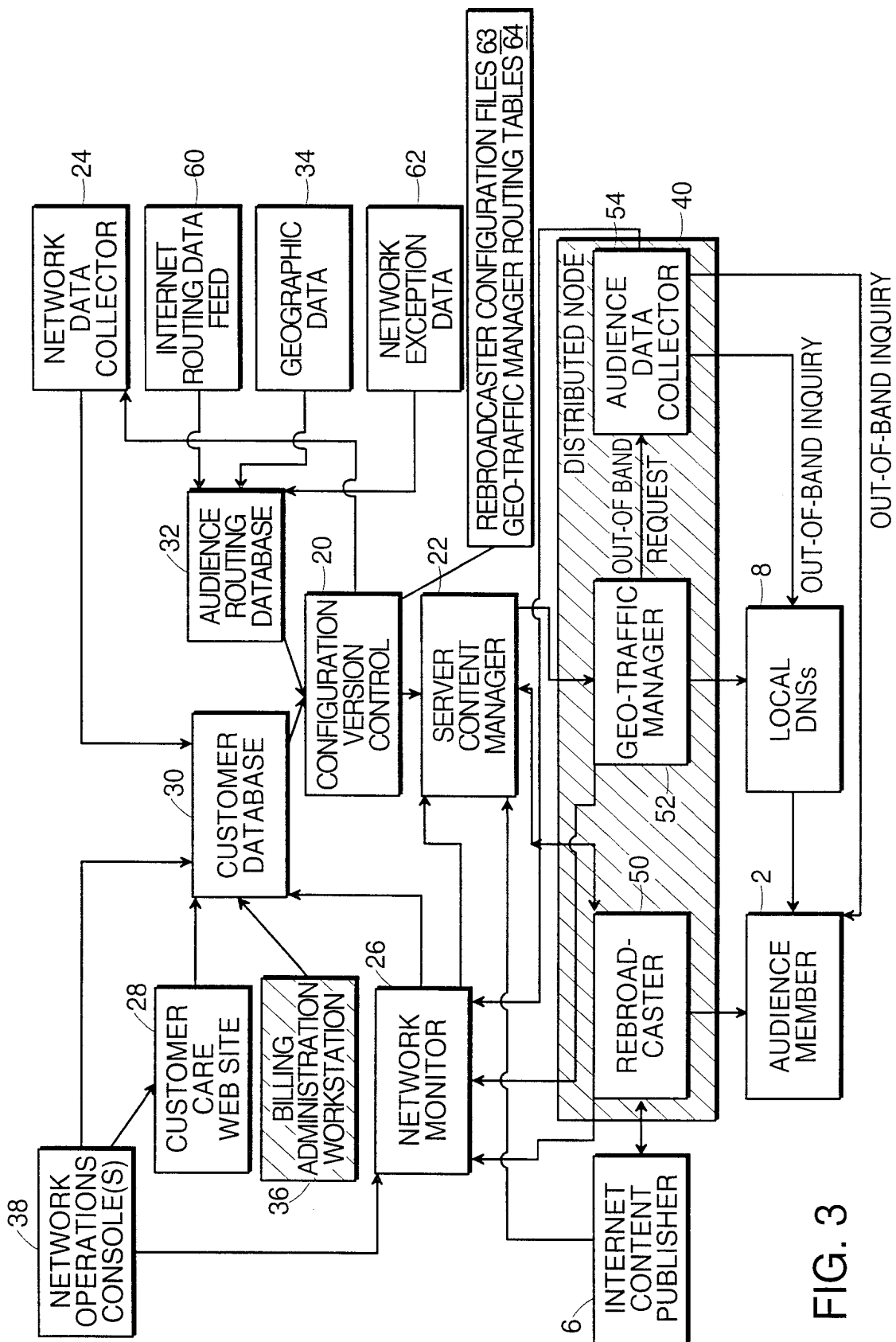


FIG. 3

4/16

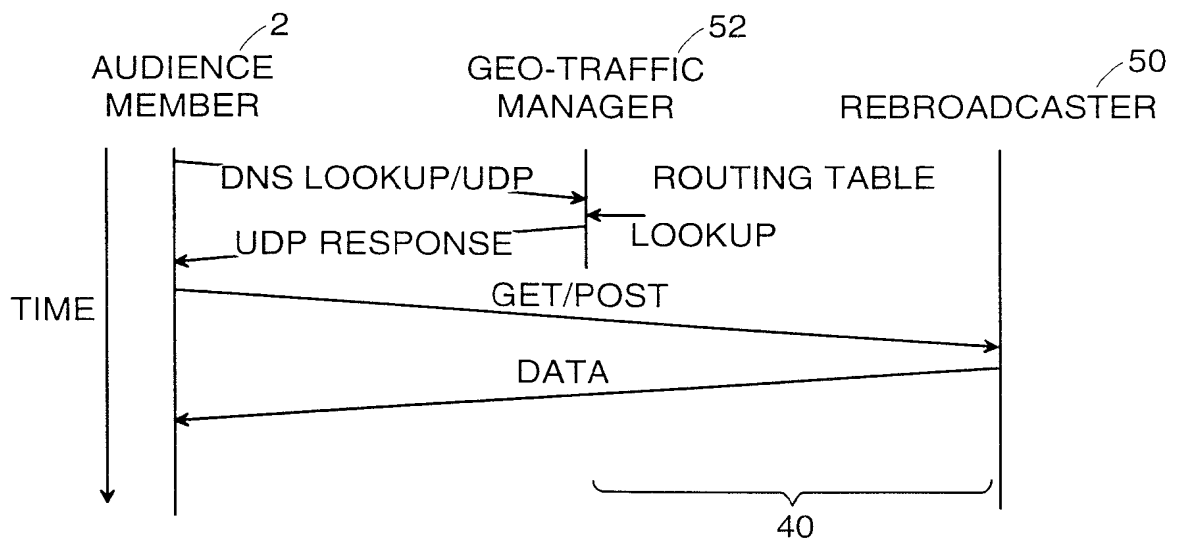


FIG. 4

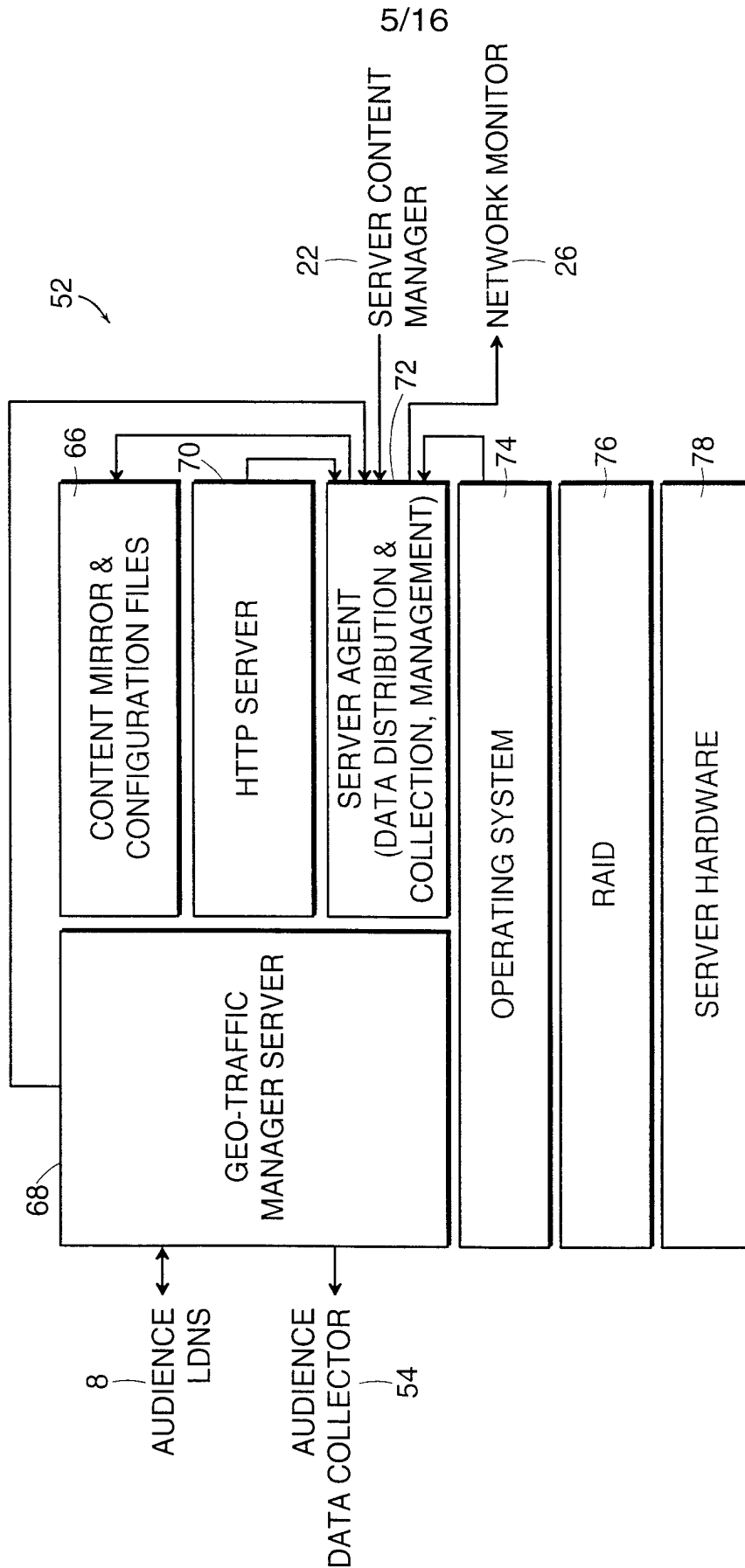


FIG. 5A

6/16

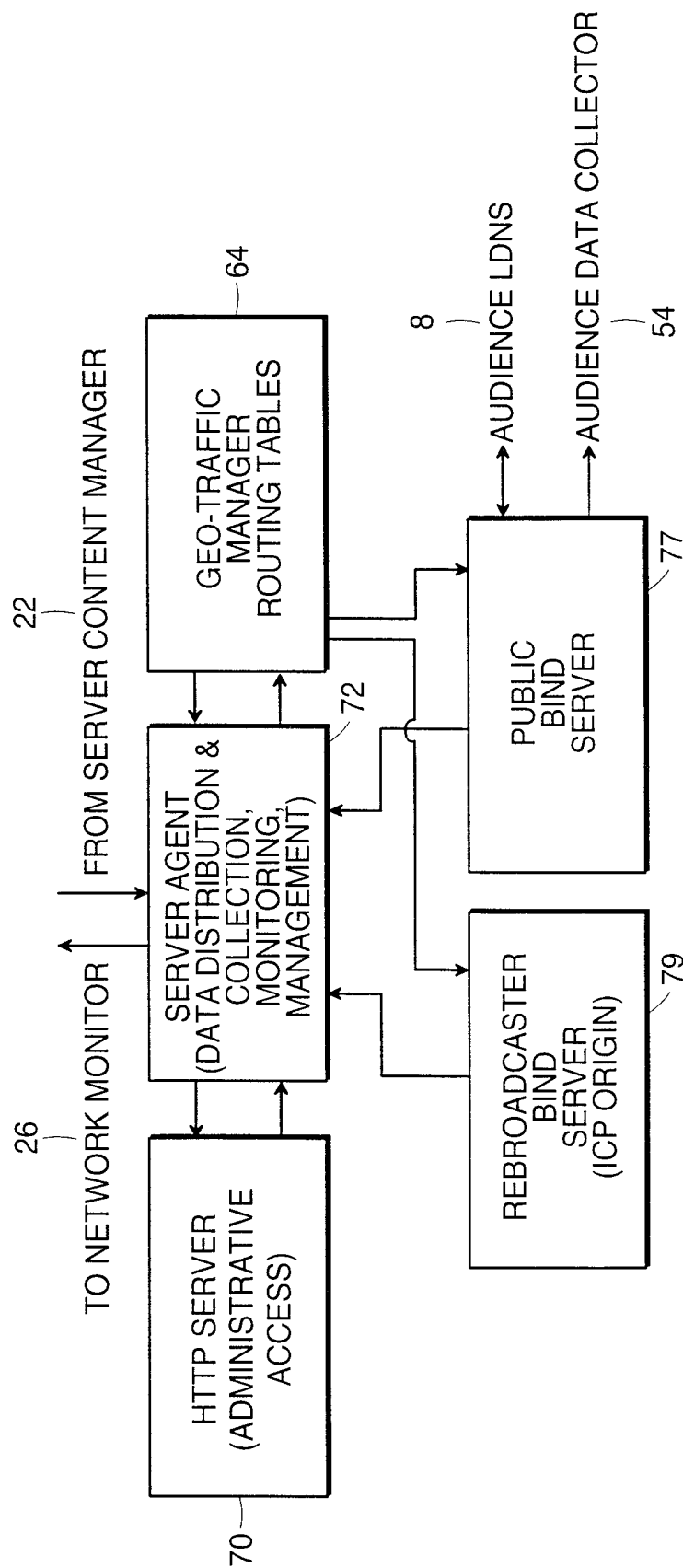


FIG. 5B

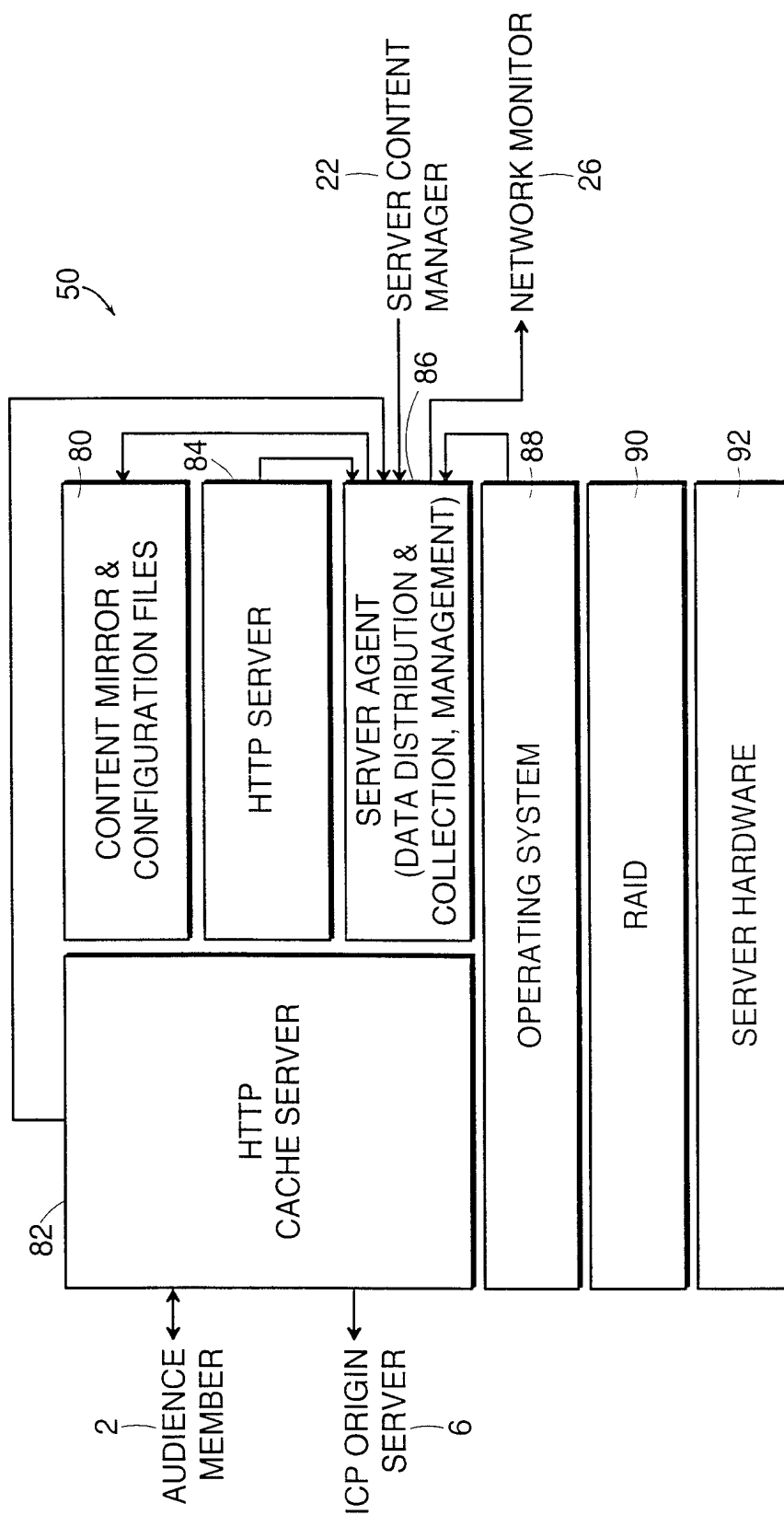


FIG. 6A

8/16

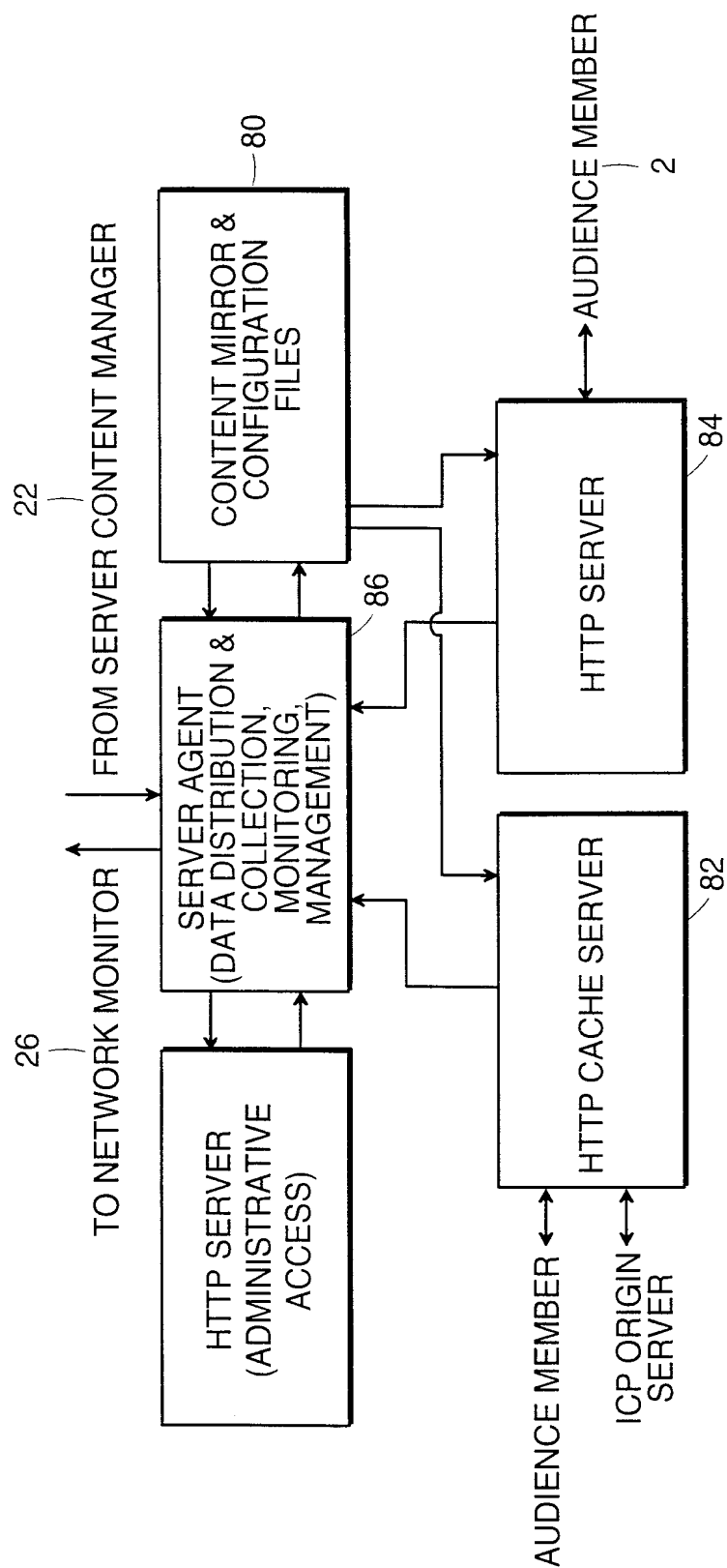


FIG. 6B

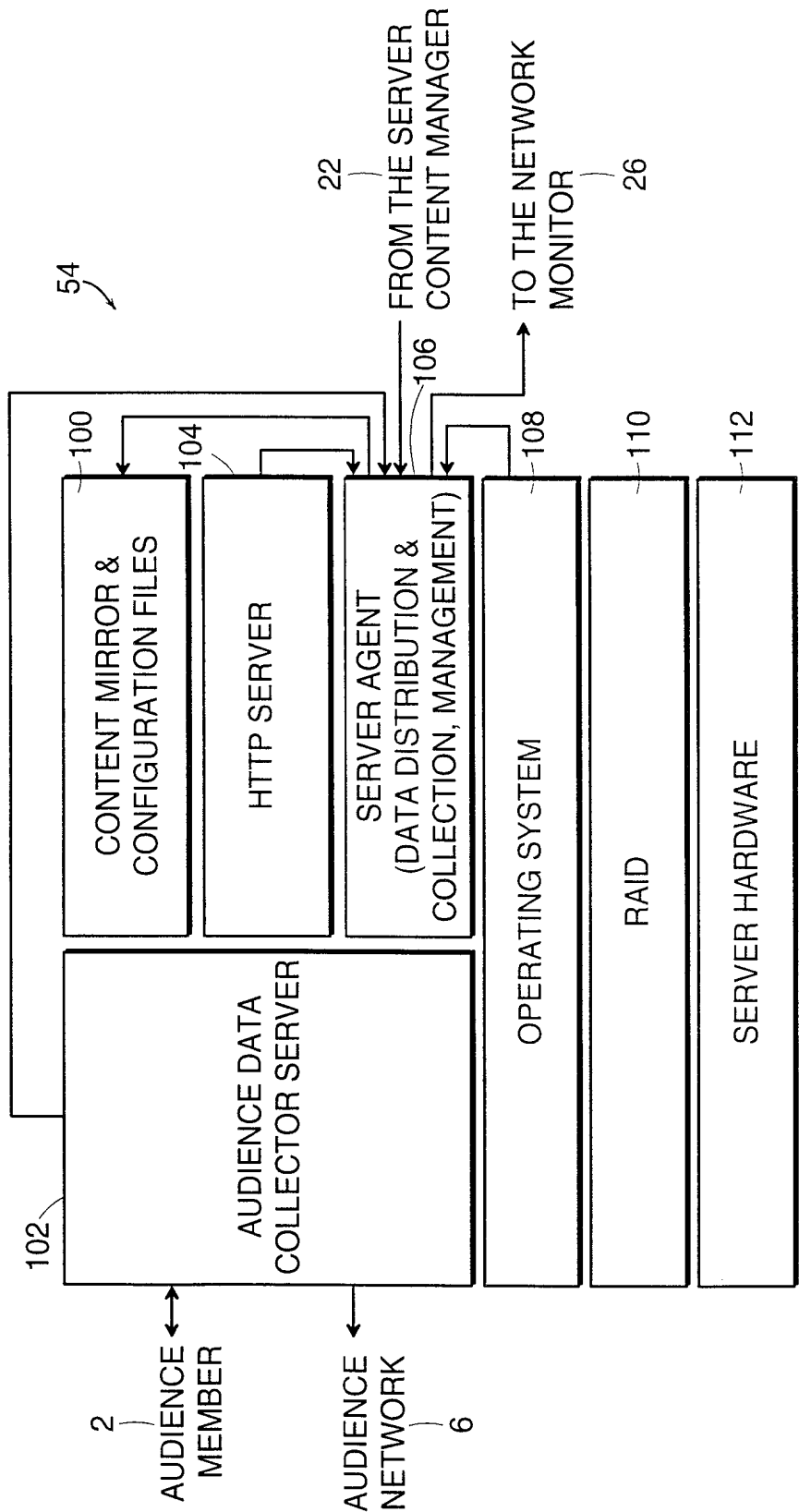


FIG. 7A

10/16

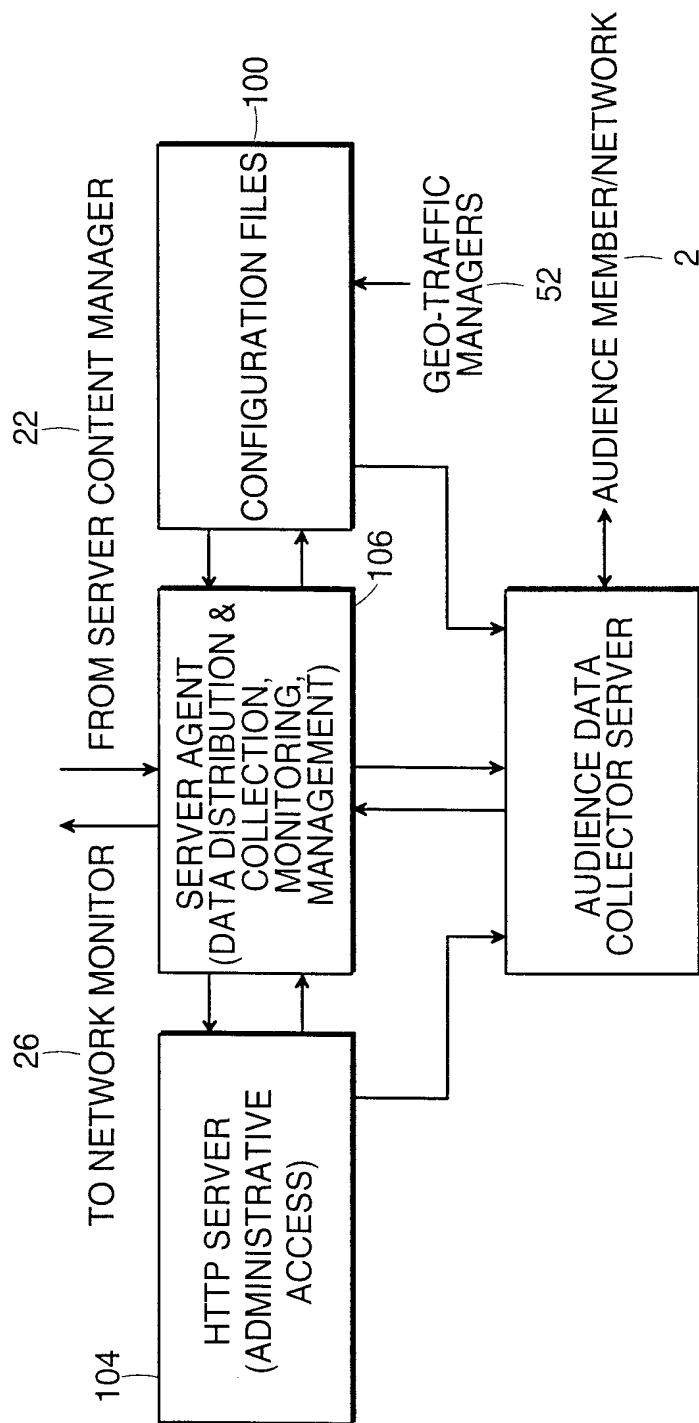


FIG. 7B

11/16

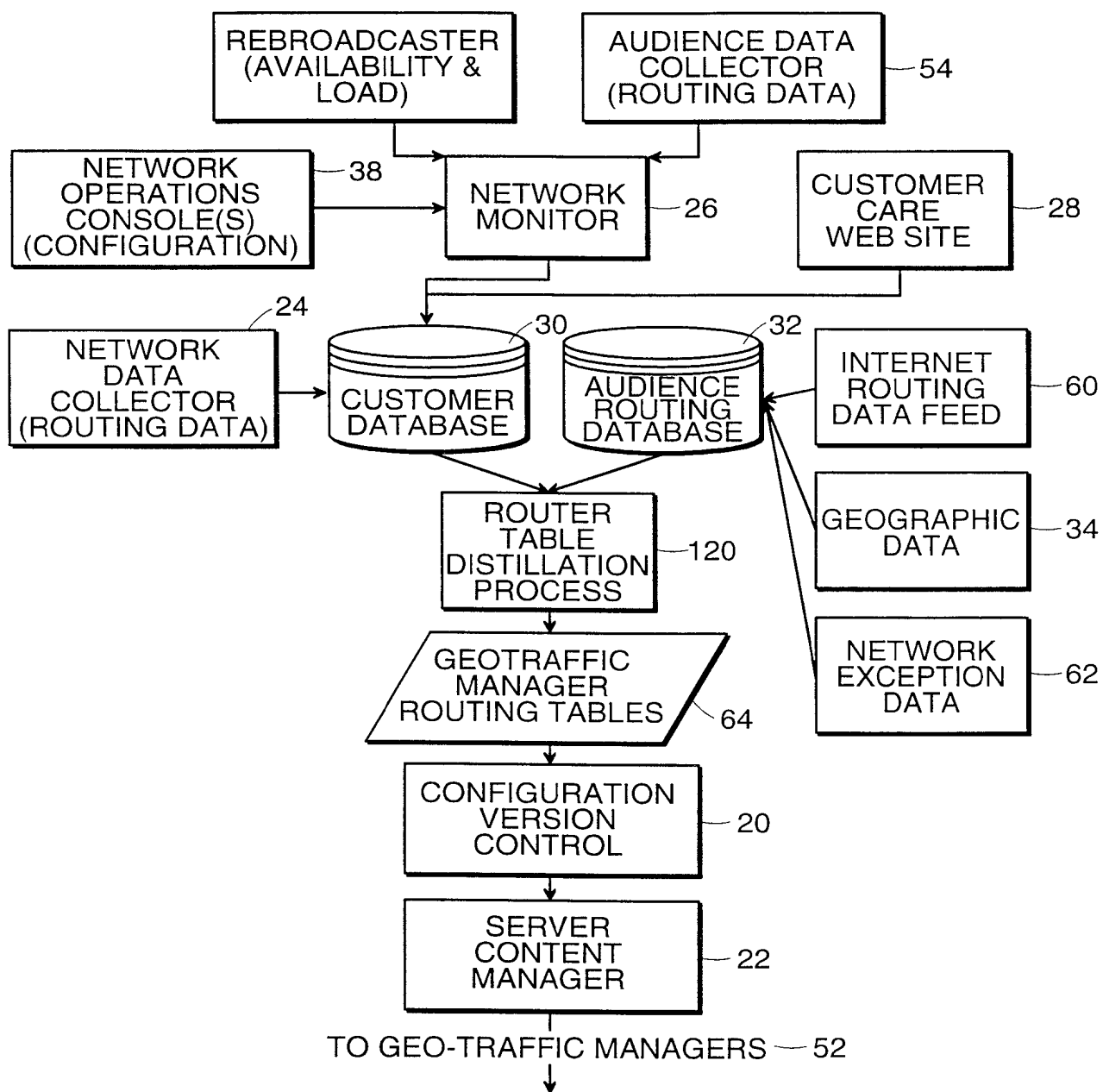


FIG. 8

12/16

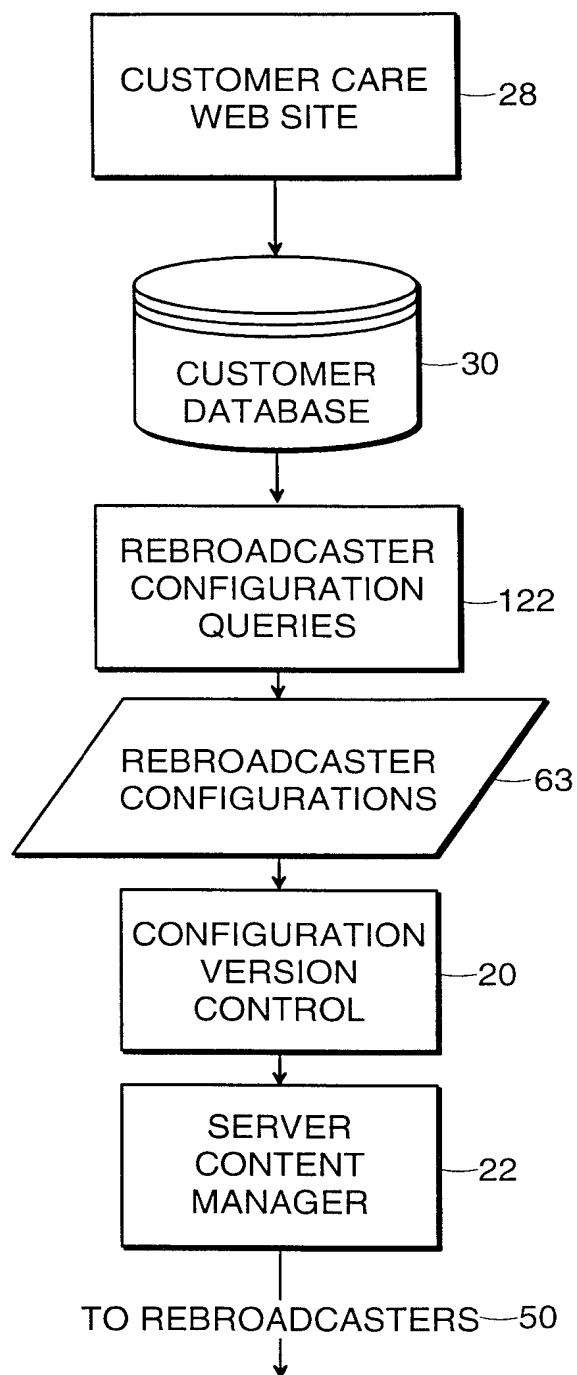


FIG. 9

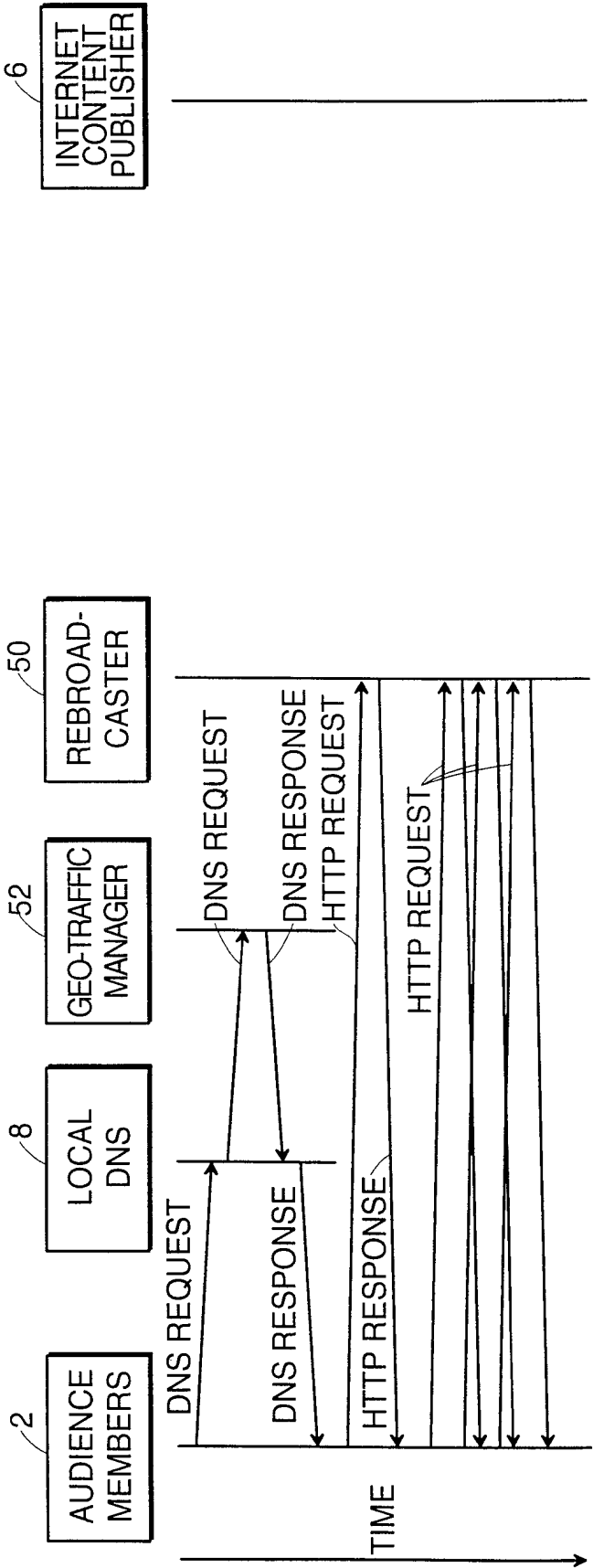


FIG. 10

14/16

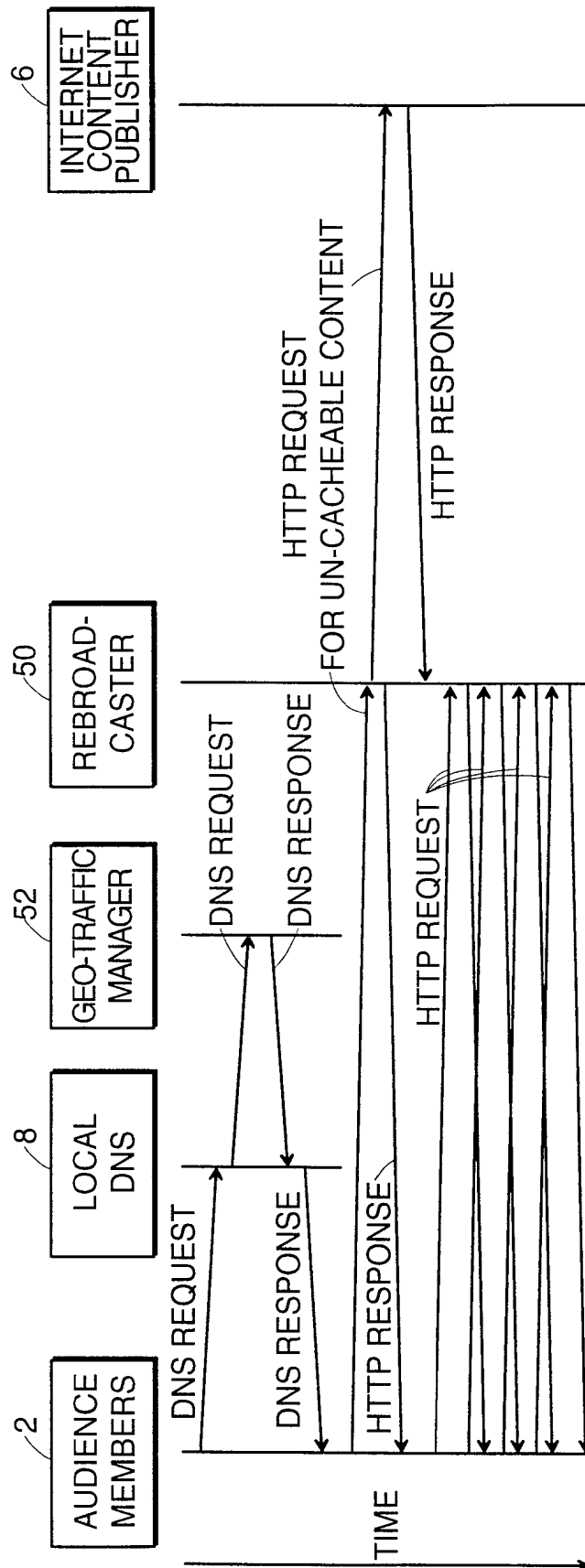


FIG. 11

15/16

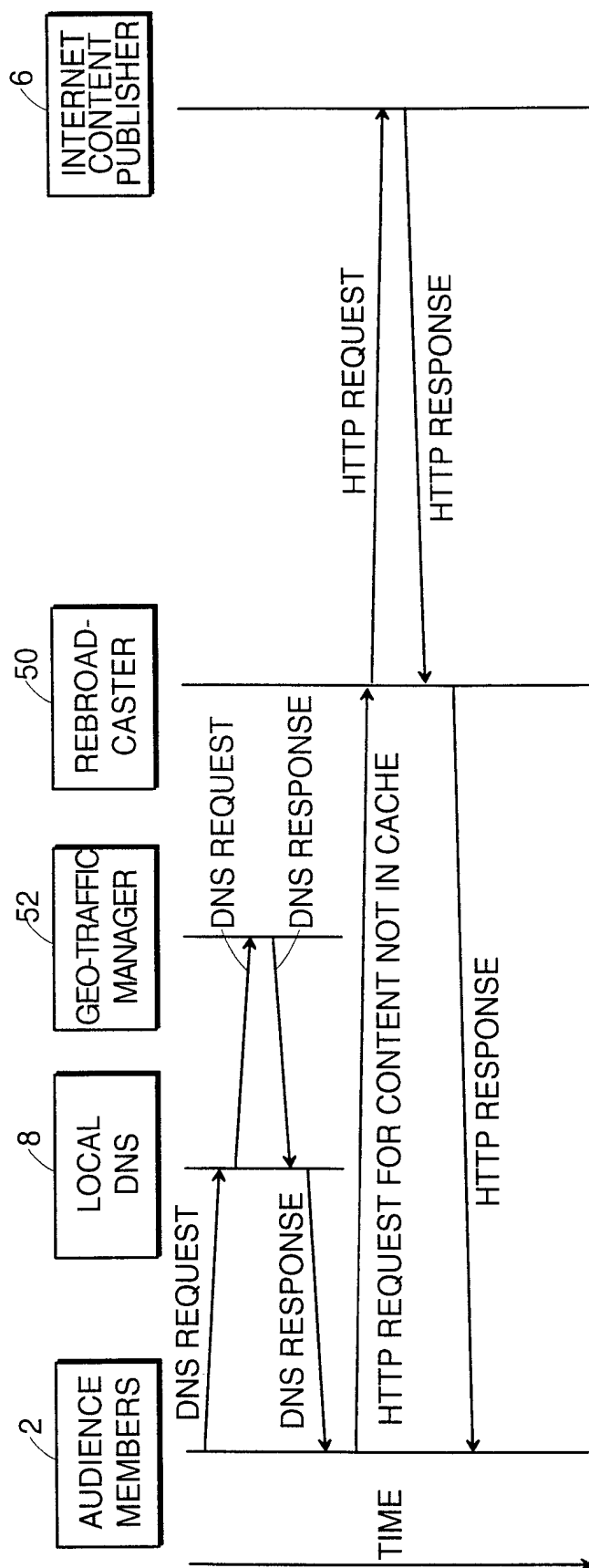


FIG. 12

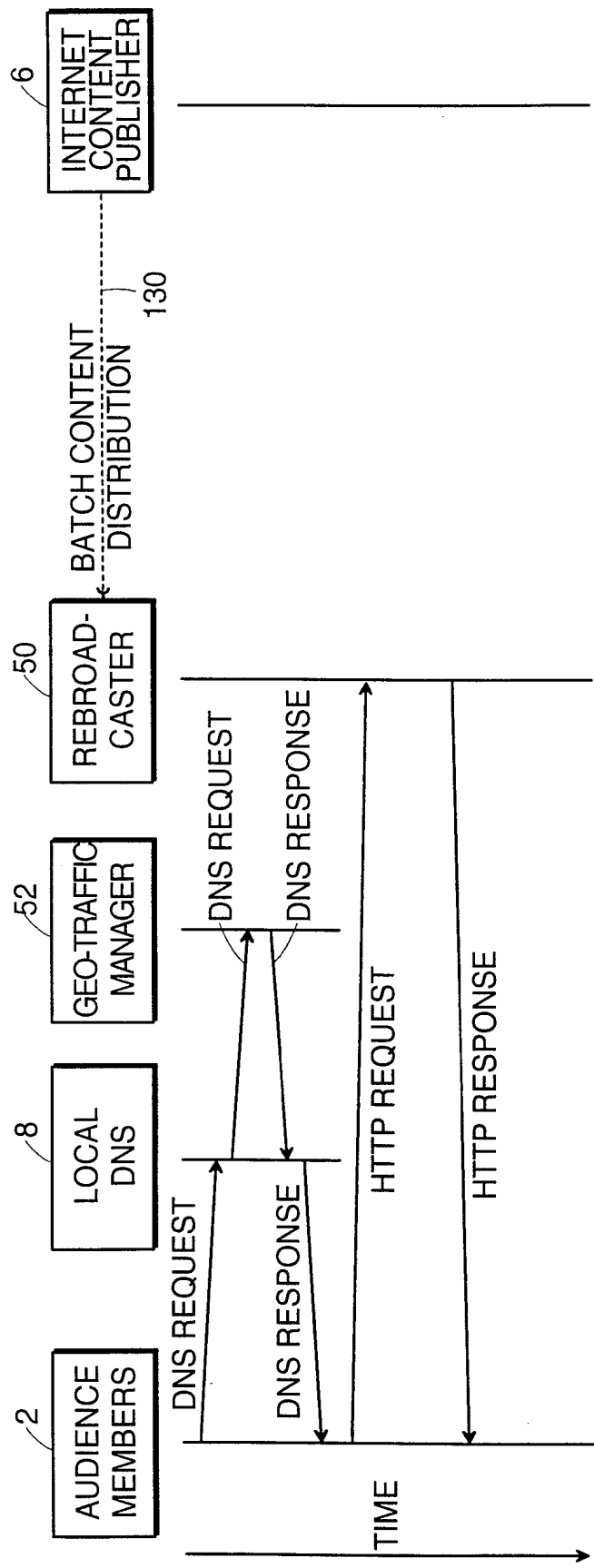


FIG. 13